

*Федеральное государственное автономное образовательное учреждение
высшего образования*
**НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»**

Образовательная программа «Прикладная математика»
бакалавр

ОТЧЕТ
по проектной работе

Распознавание пола по аудиозаписи

Выполнили студенты гр. БПМ-195

Бикушев Глеб Дмитриевич

Зубенок Павел Александрович

Шевцев Марк Денисович

Руководитель проекта:

(должность, ФИО руководителя проекта)

(оценка)

(подпись)

(дата)

Содержание

1	Вступление	3
2	Постановка задачи	3
2.1	Неформальная постановка задачи	3
2.2	Формальная постановка задачи	3
3	Методы решения задачи	3
3.1	Описание звуковых параметров	3
3.2	Рассматриваемые алгоритмы машинного обучения	4
3.3	Используемое программное обеспечение	4
4	Решение поставленной задачи	5
4.1	Создание базы данных	5
4.2	Обучение нейронной сети классификатором KNN	7
4.3	Обучение нейронной сети классификатором RFC	8
4.4	Основная программа, используемая для тестирования моделей	9
4.5	Сравнение классификаторов KNN и RFC	9
5	Тесты	10
6	Результат	11
	Список литературы	11

1 Вступление

Распознавание пола по аудиозаписи является важной частью систем автоматического распознавания речи и систем автоматического голосового ответа. Решение этой проблемы снижает вычислительную нагрузку на такие системы, позволяя обрабатывать человеческий голос быстрее и эффективнее.

Целью данной работы является изучение современных методов классификации данных, с выбором оптимальных параметров, чтобы найти оптимальный алгоритм определения пола человека по его голосу.

Предметом этого проекта является упрощение и ускорение обработки человеческой речи, маркировки данных. Создание систем распознавания человеческой речи. Используя разные методы машинного обучения, для этого в основном используются разные классы нейронных сетей. Так как для формирования нейронных сетей вручную необходимо большое количество денег и данных, используются дополнительные решения.

Этот проект должен способствовать маркировке аудиоданных по полу и используемые идеи в нем, обеспечат лучшую точность при решении проблемы перевода аудиозаписей в текст.

2 Постановка задачи

2.1 Неформальная постановка задачи

Используя эти аудиозаписи, создать обучающий набор данных, обучить модели на них и дать их сравнительные характеристики, предсказать пол человека с помощью тестовых аудиозаписей.

2.2 Формальная постановка задачи

Предоставлен набор аудиофайлов с расширением wav. Из этого набора случайным образом выберем 80% аудиозаписей и сформируем массив формата `pr.array` размером $m \times n$ - обучающий образец, где первые столбцы $n - 1$ являются значениями признаков для соответствующих фрагментов аудиозаписи, а n -й столбец пол человека чья это аудиозапись. Следовательно, m - число таких фрагментов. Полученный массив должен быть передан как обучающий набор к разным алгоритмам из библиотеки SciKit-Learn (далее - sklearn) и доказать гипотезу: "Независимо от того, на каком языке обучить модель, она будет верно определять пол диктора любого языка".

3 Методы решения задачи

3.1 Описание звуковых параметров

Мел-кепстральные коэффициенты (MFCC) были выбраны в качестве параметров звука. потому что эта количественная оценка звука основана на статистической обработке большого количества данных о восприятии высоты звука человеческим слухом. Высота звука, воспринимаемого человеческим слухом, нелинейно зависит от его частоты (см.рис 1). И хотя эта зависимость не очень точна, ее довольно удобно использовать, потому что она описывается простой формулой

$$m = 1127.01048 \ln \left(1 + \frac{f}{700} \right) \quad (1)$$

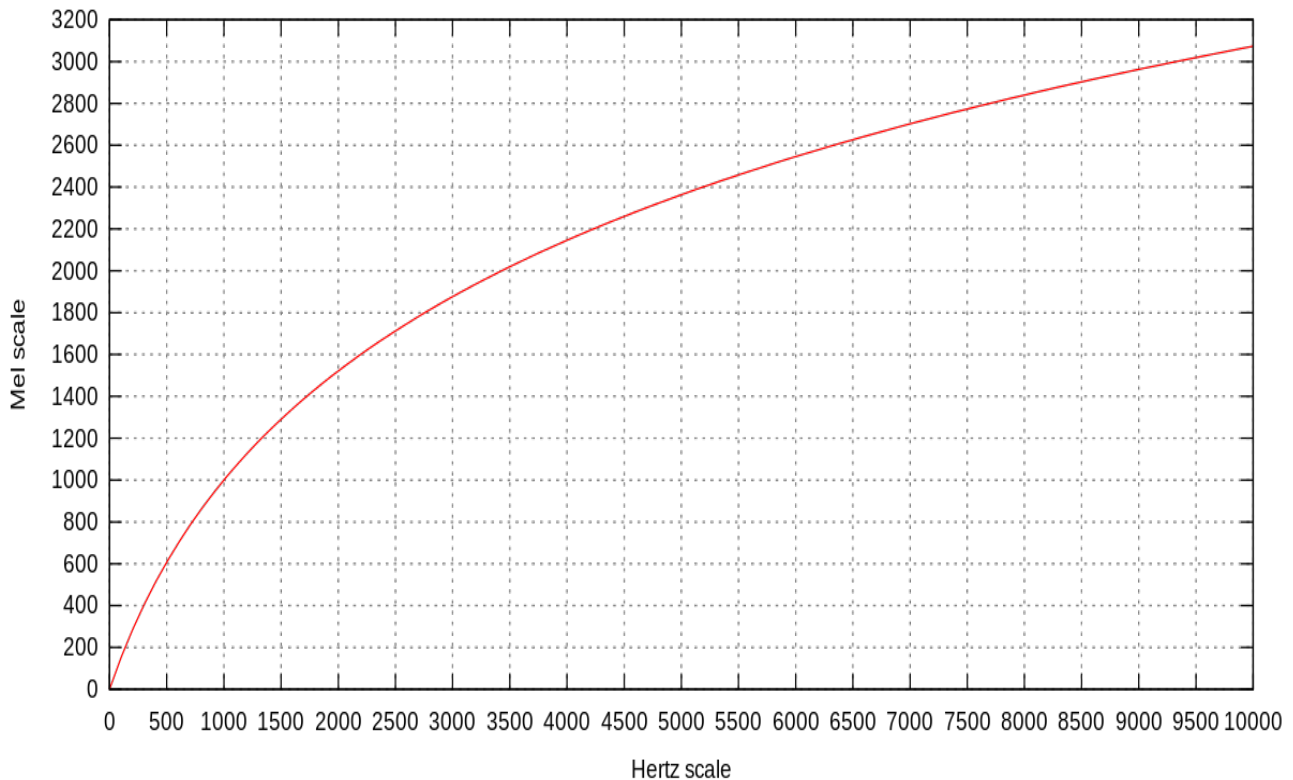


Рис. 1: График зависимости высоты звука в мелах от частоты колебаний (для чистого тона)

3.2 Рассматриваемые алгоритмы машинного обучения

Поскольку наша задача является задачей классификации, то рассматривать будем исключительно алгоритмы машинного обучения с подкреплением. Воспользуемся следующими моделями, поскольку они лучше остальных подходят под нашу задачу:

- **KNN (K Nearest Neighbors)** – метрический классификатор, основанный на оценивании сходства объектов. Классифицируемый объект относится к тому классу, которому принадлежат ближайшие к нему объекты обучающей выборки.
- **RFC (Random Forest Classifier)** – метаоценщик, который сопоставляет ряду классификаторов решающих деревьев различные подвыборки набора данных и использует усреднение для улучшения точности предсказания класса и контроля переобучения.

3.3 Используемое программное обеспечение

1. **python_speech_features** – библиотека для работы с записями голосов, используется для извлечения MEL коэффициентов из аудиозаписей.
2. **scipy.io.wavfile** – библиотека для обработки аудиозаписей формата wav. В частности, в работе из этой библиотеки была использована функция `read()`, которая получает на вход wav-файл, а возвращает частоту дискретизации (в сэмплах / сек) и данные из файла wav, которые записываются в массив NumPy
3. **Anaconda** – бесплатный дистрибутив языков Python и R для научных вычислений (анализ данных, машинное обучение и приложения, обработка данных, предиктивный анализ, и так далее).

4. **Pandas** – библиотека для работы с такими объектами, как Series и DataFrame (представляют собой таблицы данных). Объект DataFrame (таблицу) можно очень удобно перевести в формат csv и работать с отдельно с этим файлом, а так же извлекать данные из файла csv в DataFrame
5. **SciKit-Learn** – библиотека, содержащая в себе реализацию основных алгоритмов машинного обучения и проверки их на точность.
6. **Pickle** – встроенный в Python модуль, использованный для сохранения наилучших моделей по результатам обучения.
7. **Matplotlib.pyplot** – библиотека для построения необходимых по заданию графиков, визуализации данных.
8. **NumPy** – библиотека, упрощающая и ускоряющая работу с вычислениями. Также необходим для полноценной работы с остальными модулями, позволяет удобно хранить массивы большой базы данных и обращаться к ним
9. **speech_recognition** и **pyaudio** – библиотеки, необходимые для записи голоса непосредственно в программе

4 Решение поставленной задачи

4.1 Создание базы данных

Для реализации машинного обучения (в нашем случае речь идёт об обучении "с учителем" поскольку мы используем классификаторы) необходима база помеченных данных, в которой важная часть признаков уже разделена на отдельные категории или классы. Поэтому создадим такую базу данных, которая состоит из извлеченных MEL коэффициентов из аудиозаписей. Возьмём 16 аудиозаписей (8 - мужских, 8 - женских). Прежде всего, поясним, какие параметры функции `mfcc()` были взяты и для чего.

- **sig** – звуковой сигнал для расчета функций.
- **rate** – частота дискретизации.
- **winlen** - длина окна анализа в секундах. По умолчанию 0,025 с (25 миллисекунд). Сделаем этот показатель на уровне 0.2 с.
- **winstep** - шаг между последовательными окнами в секундах. По умолчанию 0,01 с (10 миллисекунд). Положим это значение, равное 0.01.
- **nfilt** - количество фильтров в наборе фильтров, по умолчанию 26. Для нашей задачи подходит 26 фильтров. Оставим это значение.
- **nfft** - размер БПФ. По умолчанию 512. Нам понадобится значение 12512.
- **lowfreq** - нижняя граница полосы фильтров mel. В Гц по умолчанию 0. Такое значение нам подходит.
- **highfreq** - максимальный край полосы фильтров из mel. Для нашей задачи подойдет значение 22050 Гц.
- **preemph** - применить фильтр preemphasis с preemph в качестве коэффициента. По умолчанию 0,97. Нам не нужен данный фильтр. Выставим значение 0.

Код программы:

```
1  from python_speech_features import mfcc
2  import scipy.io.wavfile as wav
3  import numpy as np
4  import pandas as pd
5  from sklearn import preprocessing
6
7  def make_data(name, dfa):
8      (rate, sig) = wav.read(name)
9      mfcc_feat = mfcc(sig, rate, winlen=0.2, winstep=0.01,
10                      nfilt=26, nfft=12512, lowfreq=0,
11                      highfreq=22050, preemph=0)
12
13      res_shape = np.shape(mfcc_feat)
14      df = pd.DataFrame(mfcc_feat)
15      if name.find("woman") != -1:
16          df['SEX'] = pd.Series([0] * res_shape[0], index=df.index)
17      elif name.find("man") != -1:
18          df['SEX'] = pd.Series([1] * res_shape[0], index=df.index)
19      else:
20          print("Wrong name of the recording. ")
21          exit(1)
22      dfa = dfa.append(df)
23      return dfa
24
25  # Создаем пустой датафрейм, который и будем передавать функциям
26  df_all = pd.DataFrame()
27
28  i = 1
29  while i < 9:
30      df_all = make_data("./training/man.wav".format(i), df_all)
31      df_all = make_data("./training/woman.wav".format(i), df_all)
32      i += 1
33
34  # переводим полученный датафрейм в формат csv
35  df_all.to_csv("filedata.csv")
36
37  # нормируем данные файла filedata.csv
38  min_max_scaler = preprocessing.MinMaxScaler()
39  np_scaled = min_max_scaler.fit_transform(df_all)
40  df_normalized = pd.DataFrame(np_scaled)
41  df_normalized.to_csv("filedata_norm.csv")
```

На выходе получаем два файла формата csv : "filedata.csv" и "filedata_norm.csv" ("filedata.csv" , данные которого отнормированы)

4.2 Обучение нейронной сети классификатором KNN

Код программы:

```
1      import numpy as np
2      import pandas as pd
3      import pickle
4      from sklearn.model_selection import GridSearchCV
5      from sklearn.neighbors import KNeighborsClassifier
6      from sklearn.metrics import precision_score
7
8      frame = pd.read_csv("filedata_norm.csv", delimiter=',', index_col=0)
9
10     x_train = frame.iloc[:, :13].to_numpy()
11     y_train = frame.iloc[:, 13].to_numpy()
12
13     tuned_parameters = {
14         "weights": np.array(["uniform " "distance"]),
15         "algorithm": np.array(["ball_tree " "kd_tree " "brute"]),
16         "p": np.array([2, 12]),
17     }
18
19     model = GridSearchCV(KNeighborsClassifier(), tuned_parameters, cv=14)
20     # cv – количество блоков кросс-валидации
21     model.fit(x_train, y_train)
22
23     best_model = model.best_estimator_
24     y_true = y_train
25     y_pred = best_model.predict(x_train)
26     precision = precision_score(y_true, y_pred, average='weighted')
27     print(precision)
28     if precision >= 0.95:
29         with open("knn_params.pickle", "wb") as fout:
30             pickle.dump(best_model.get_params(), fout)
31             fout.close()
32         with open("knn_best_model.sav", "wb") as fout:
33             pickle.dump(best_model, fout)
34             fout.close()
```

После компиляции создаются 2 файла : "knn_params.pickle" и "knn_best_model.sav", которые хранят информацию о наилучшей модели в результате обучения.

4.3 Обучение нейронной сети классификатором RFC

Код программы:

```
1      import pandas as pd
2      import numpy as np
3      import pickle
4      from sklearn.ensemble import RandomForestClassifier
5      from sklearn.model_selection import GridSearchCV
6      from sklearn.metrics import precision_score
7
8      frame = pd.read_csv("filedata_norm.csv" , delimiter=',', index_col=0)
9
10     x_train = frame.iloc[:, :13].to_numpy()
11     y_train = frame.iloc[:, 13].to_numpy()
12
13     tuned_parametres = {
14         "n_estimators": np.array([50, 100, 150]),
15         "criterion": np.array(["gini "entropy"]),
16         "max_features": np.array(["sqrt "log2 None])
17     }
18
19     model = GridSearchCV(RandomForestClassifier(), tuned_parametres, cv=14)
20
21     model.fit(x_train, y_train)
22     best_model = model.best_estimator_
23     y_true = y_train
24     y_pred = best_model.predict(x_train)
25     precision = precision_score(y_true, y_pred, average='weighted')
26     print(precision)
27
28     if precision >= 0.95:
29         with open("rfc_params.pickle" , "wb" ) as fout:
30             pickle.dump(best_model.get_params(), fout)
31             fout.close()
32         with open("rfc_best_model.sav" , "wb" ) as fout:
33             pickle.dump(best_model, fout)
34             fout.close()
```

После компиляции создаются 2 файла : "rfc_params.pickle" и "rfc_best_model.sav" , которые хранят информацию о наилучшей модели в результате обучения.

4.4 Основная программа, используемая для тестирования моделей

Она состоит из функции `make_data`, которая используется для извлечения из аудиозаписи `mel` - коэффициентов.

Далее идёт функция `voice_to_numpy`, которая нормирует `mel` - коэффициенты и записывает получившийся результат в массив `NumPy` :

```
1 def voice_to_numpy(name):
2     frame = pd.DataFrame()
3     frame = make_data(name, frame)
4     min_max_scaler = preprocessing.MinMaxScaler()
5     np_scaled = min_max_scaler.fit_transform(frame)
6     frame_normalized = pd.DataFrame(np_scaled)
7     data = frame_normalized.to_numpy()
8     return data
9
```

Затем выбрав, один из двух классификаторов извлекаем содержимое лучшей модели формата `pickle`:

```
1 if key == 1:
2     with open("knn_best_model.sav" , "rb" ) as fin:
3         model = pickle.load(fin)
4 elif key == 2:
5     with open("rfc_best_model.sav" , "rb" ) as fin:
6         model = pickle.load(fin)
7
```

В завершении считаем среднее арифметическое по столбцам предсказания пола с помощью функции библиотеки `NumPy` "`mean()`" :

```
1 result = np.mean(model.predict(data))
2 if result < 0.5:
3     print(">Женщина" )
4     print(result)
5 elif result > 0.5:
6     print(">Мужчина" )
7     print(result)
8 else:
9     print("Неопределенно." )
    print(result)
```

4.5 Сравнение классификаторов KNN и RFC

Входные данные (аудиозапись русских дикторов монозвука с расширением `wav`) представляют собой временные ряды - последовательность измерений в течение каждого измерения, характеризуется соответствующим вектором коэффициентов `MEL`. Рассмотрим

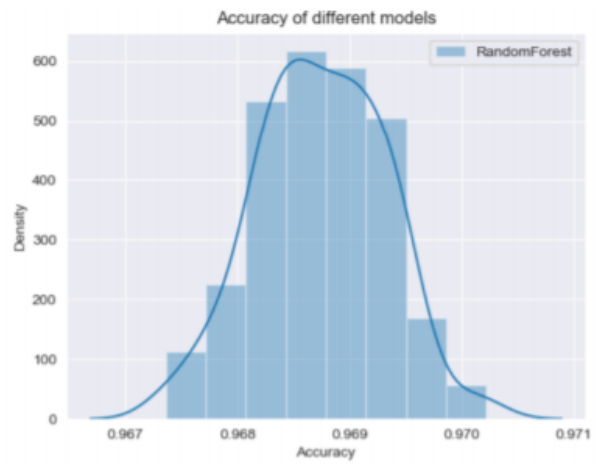
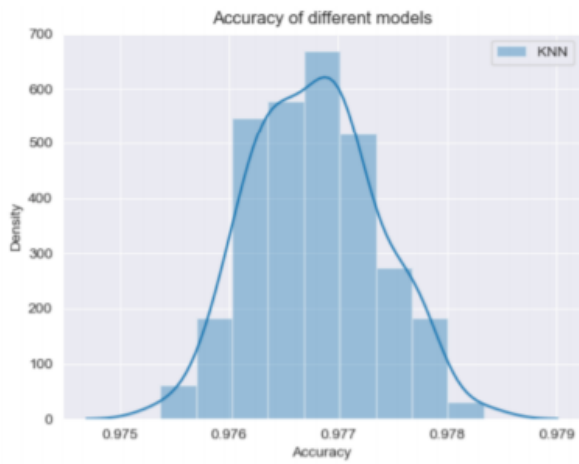


Рис. 2: Сравнительные графики зависимости плотности вероятности точности модели от её значения (за 100 запусков)

эффективность алгоритмов машинного обучения и построим сравнительные графики по стандартным параметрам с помощью библиотеки Matplotlib.

Из рис. 2 можно наблюдать, что наши модели KNN и RFC показывают точность порядка 97%, что является отличным показателем для машинного обучения.

5 Тесты

В предыдущем пункте мы сравнивали 2 алгоритма (knn и rfc). Входные данные были аудиозаписи на русском языке.

Ниже представлена таблица где входные данные были аудиозаписи на разных языках.

Как можно заметить, хоть и с небольшой погрешностью, но разработанная модель верно определяет пол диктора иностранного языка, не смотря на то, что была обучена на записях русской речи.

номер записи	язык диктора	классификатор	пол	результат
1	английский	KNN	муж	жен (0.243)
1	английский	RFC	муж	жен (0.342)
2	английский	KNN	жен	жен (0.267)
2	английский	RFC	жен	жен (0.187)
3	немецкий	KNN	муж	муж (0.509)
3	немецкий	RFC	муж	муж (0.564)
4	немецкий	KNN	жен	жен (0.438)
4	немецкий	RFC	жен	жен (0.213)
5	французский	KNN	муж	муж (0.611)
5	французский	RFC	муж	муж (0.506)
6	французский	KNN	жен	жен (0.335)
6	французский	RFC	жен	жен (0.212)
7	испанский	KNN	муж	муж (0.619)
7	испанский	RFC	муж	муж (0.758)
8	испанский	KNN	жен	жен (0.473)
8	испанский	RFC	жен	жен (0.292)
9	вьетнамский	KNN	муж	муж (0.561)
9	вьетнамский	RFC	муж	муж (0.230)
10	вьетнамский	KNN	жен	жен (0.387)
10	вьетнамский	RFC	жен	жен (0.244)
11	грузинский	KNN	муж	муж (0.681)
11	грузинский	RFC	муж	муж (0.565)
12	грузинский	KNN	жен	жен (0.438)
12	грузинский	RFC	жен	муж (0.723)
13	айзербайджанский	KNN	муж	муж (0.554)
13	айзербайджанский	RFC	муж	муж (0.703)
14	айзербайджанский	KNN	жен	жен (0.476)
14	айзербайджанский	RFC	жен	жен (0.408)

6 Результат

Была разработана программа, которая определяет пол диктора по аудиозаписи. Была проверена и подтверждена гипотеза: "Независимо от того, на каком языке обучить модель, она будет верно определять пол диктора любого языка" . Найти код проекта, а также самому проверить исправность работы нашей модели, записав свой голос сразу в программе можно по этой [ссылке](#), перейдя к GitHub репозиторий.

Список литературы

- [1] Мел-кепстральные коэффициенты (MFCC) и распознавание речи // article from Habr.com // URL: <https://habr.com/ru/post/140828/>
- [2] Классификатор kNN // article from Habr.com.// URL: <https://habr.com/ru/post/149693/>
- [3] Реализация и разбор алгоритма «случайный лес» на Python // article from tproger.ru// URL: <https://tproger.ru/translations/python-random-forest-implementation/>
- [4] Welcome to python_speech_features's documentation! // URL: <https://python-speech-features.readthedocs.io/en/latest/>